# Optimization of CUDA GPU Kernels and Translation to AMDGPU in Polygeist/MLIR

Ivan R. Ivanov[1,2]

Alex Zinenko[3], Jens Domke[2], Endo Toshio[1], William S. Moses[4]

1. Tokyo Institute Of Technology
2. RIKEN R-CCS
3. Google
4. UIUC

# Motivation

- Writing high performance CUDA code is **hard**
- Even more difficult to make it portable
  Big differences in GPU archs - warp size is 16, 32, 64 on Intel,NVIDIA, AMD

| GPU | Consumer-grade | | HPC | |
|---|---|---|---|---|
| | NVIDIA A4000 | AMD RX6800 | NVIDIA A100 | AMD MI210 |
| Compute Capability | 8.6 | gfx1030 | 8.0 | gfx90a |
| SMs | 48 | 60 | 108 | 104 |
| FLOPs (f64) | 0.60T | 1.01T | 9.75T | 22.60T |
| FLOPs (f32) | 19.17T | 16.17T | 19.49 | 22.60T |
| Memory Bandwidth | 445 GB/s | 512 GB/s | 1555 GB/s | 1638 GB/s |
| Global Memory | 16 GB | 16 GB | 40 GB | 64 GB |
| L2 Cache | 4 MB | 4MB | 40 MB | 16 MB |
| L1 Cache (Per SM) | 128 KB | 16 KB | 192 KB | 16 KB |

- Large amount of (legacy) scientific C/C++ CUDA code
- **Large cost** of porting and tuning to another GPU architecture (or vendor)

# Aim and Contributions

**Write simple CUDA code once** - let the compiler optimize and tune it for the target architecture (even AMD GPUs!)
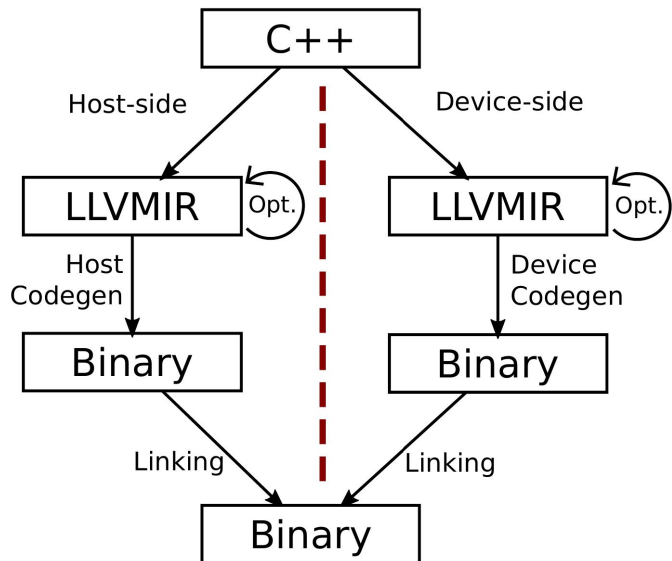
How?

- **Represent parallel GPU computation in MLIR** (Extend Polygeist)
- Optimizations for **maximizing target hardware utilization**
- **CUDA to AMDGPU translation**

# Traditional compilers

Device and Host side code split at the start of the pipeline (Example: clang)
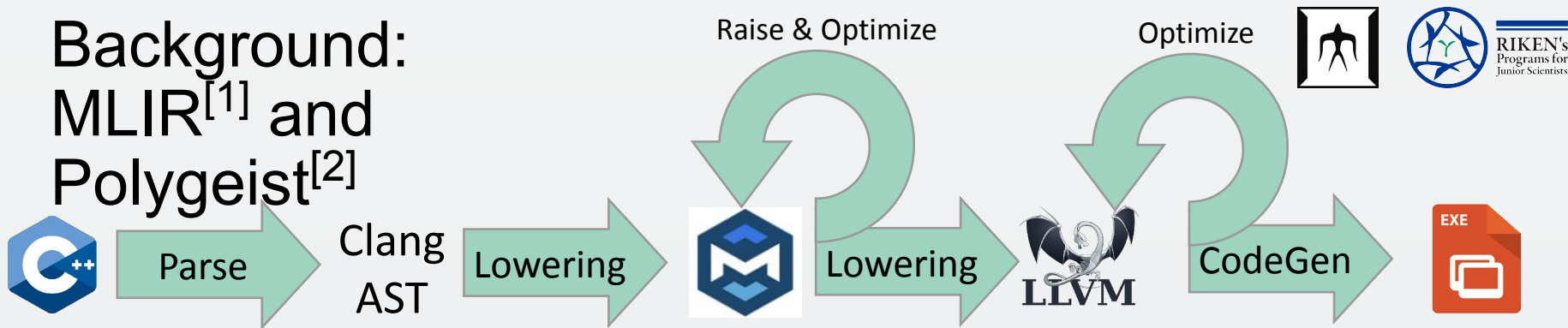


The information about the parallelism is hidden behind a library call

```
; Host side
declare @foo
define void @main() {
  ...
  call @cudaLaunchKernel(... , @foo)
}
```

```
; Device side
define void @foo(%a, %b, %c, %d, %i) {
  %s = call @sqrt(%d) ; Executed N times
  c[%i] = a[%i] + b[%i] + %s
}
```

The single thread work, where is the parallelism?

# Background: MLIR[1] and Polygeist[2]



Parse → Clang AST → Lowering → (MLIR) Raise & Optimize → Lowering → LLVM Optimize → CodeGen → EXE

- Generic C and C++ frontend that generates "standard" and user-defined MLIR (templates, classes, unions, etc. all supported)

- Preserves the structure of programs (parallelism, control flow, etc)

- Collection of high-level optimization and analysis passes

[1] MLIR: Scaling Compiler Infrastructure for Domain Specific Computation. CGO'21.
[2] Polygeist: Raising C to Polyhedral MLIR; Moses, Chelini, Zhao, and Zinenko. PACT '21.

# Optimization Friendly Parallel Representation in MLIR

GPU computation with
**parallel semantics**
and **context**

- Easier parallel optimizations
- Host-device
  cross-optimizations
- Tweaking kernel launch
  configuration

[1] High-Performance GPU-to-CPU Transpilation and Optimization via High-Level Parallel Constructs. William S. Moses, Ivan R. Ivanov, Jens Domke, Toshio Endo, Johannes Doerfert, and Oleksandr Zinenko. (PPoPP '23).

device region

grid

shared memory

block

synchronisation

```
1  func @launch(%h_out : memref<?xf32>,
2               %h_in  : memref<?xf32>, %n : i64) {
3    // Host code
4    %d_out = gpu.alloc ...
5    %d_in = gpu.alloc ...
6    gpu.memcpy %d_in %h_in
7    polygeist.gpu_region {
8      // Device code
9      parallel.for (%bx, %by, %bz) = (0, 0, 0)
10             to (grid.x, grid.y, grid.z) {
11       %shared_val = memref.alloca : memref<f32>
12       parallel.for (%tx, %ty, %tz) = (0, 0, 0)
13               to (blk.x, blk.y, blk.z) {
14         if %tx == 0 {
15           %sum = func.call @sum(%d_in, %n)
16           memref.store %sum, %shared_val[] : memref<f32>
17         }
18         polygeist.barrier(%tx, %ty, %tz)
19         %tid = %bx + blk.x * %tx
20         if %tid < %n {
21           %res = ...
22           memref.store %res, %d_out[%tid] : memref<?xf32>
23         }
24       }
25     }
26   }
27   gpu.memcpy %h_out %d_out
28 }
```

# Support for GPU Compilation in Polygeist: The pipeline
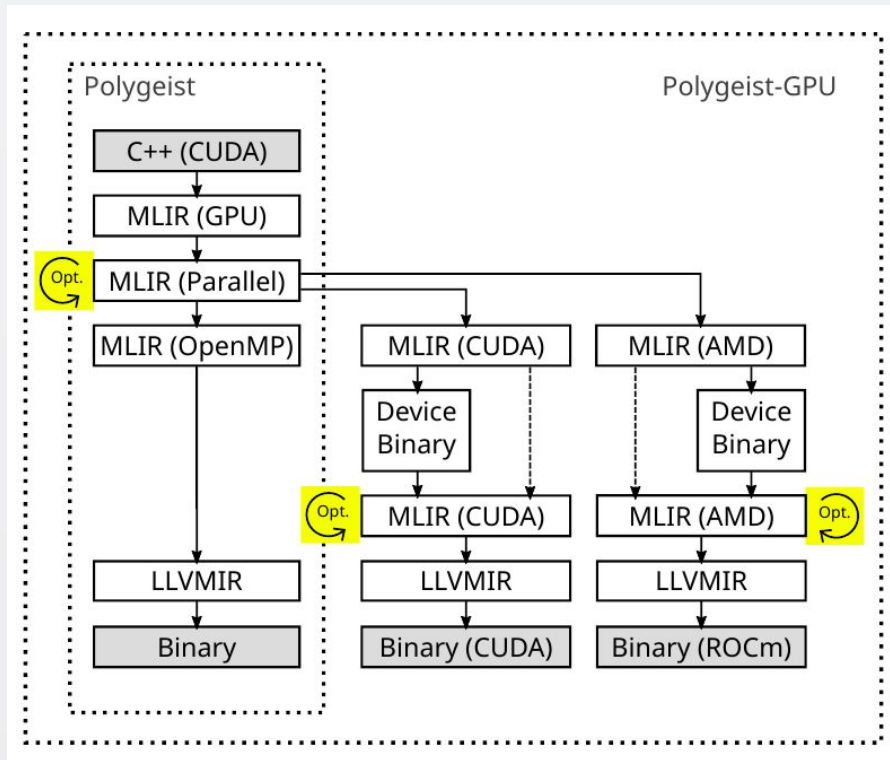
### Added

- Parallel -> Generic GPU

### Adapted from upstream MLIR

- Generic GPU -> CUDA
- Generic GPU -> AMDGPU

# Thread coarsening

Combining multiple threads worth of work in a single one (and interleaving)

```
1  parallel %i = from 0 to 16 {
2    A(%i)
3    B(%i)
4    C(%i)
5  }
```

⇨

```
1   parallel %i = from 0 to 8 {
2     %i_0 = %i * 2
3     %i_1 = %i * 2 + 1
4     A(%i_0)
5     A(%i_1)
6     B(%i_0)
7     B(%i_1)
8     C(%i_0)
9     C(%i_1)
10  }
```

+ Better instruction level parallelism, result re-use, etc..
– Worse memory access patterns…

8

# Recursive Unroll and Interleave

Generic recursive thread coarsening that works on parallel loops

```
1  parallel %i = from 0 to 16 {
2    ...
3    for %j = from 0 to 32 {
4      A(%i, %j)
5      B(%i, %j)
6    }
7  }
```

⇨

```
1  parallel %i = from 0 to 8 {
2    ...
3    for %j = from 0 to 32 {
4      A(%i_0, %j)
5      A(%i_1, %j)
6      B(%i_0, %j)
7      B(%i_1, %j)
8    }
9  }
```

# Recursive Unroll Interleave

We have a generic version of thread coarsening…
We can apply it to blocks!

```
1   parallel %block = 0 to %n {
2     %shmem = alloca()
3     parallel %thread = 0 to 1024 {
4       A(%block, %thread)
5       B(%block, %thread)
6     }
7   }
```

```
1    parallel %block = 0 to (%n / 2) {
2      %block_0 = %block * 2
3      %block_1 = %block * 2 + 1
4      %shmem_0 = alloca()
5      %shmem_1 = alloca()
6      parallel %thread = 0 to 1024 {
7        A(%block_0, %thread)
8        A(%block_1, %thread)
9        B(%block_0, %thread)
10       B(%block_1, %thread)
11     }
12   }
```

+ Nicer memory access patterns
+ **Finer grained** control over the **unroll factor** (no need to be a divisor of the bound)
   we instead generate and epilogue loop (a new kernel)

# Recursive Unroll Interleave

We have a generic version of thread coarsening…
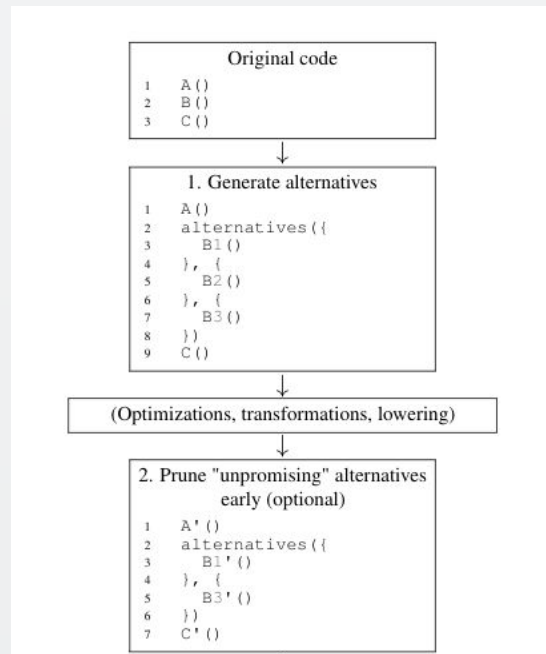
We can apply it to blocks!
We can apply it to threads!

We can **combine** the two approaches

How do we choose the best configuration from all the possible combinations?

# Alternative Code Paths



- **Multiple versions** of code that achieve the **same result** (different {block, thread} coarsening) in the parallel representation
- Defer choosing best one until later in the pipeline
  - Lower down to gpu binaries
  - Gather statistics about the kernel (registers used, memory spilled, theoretical occupancy, etc.)
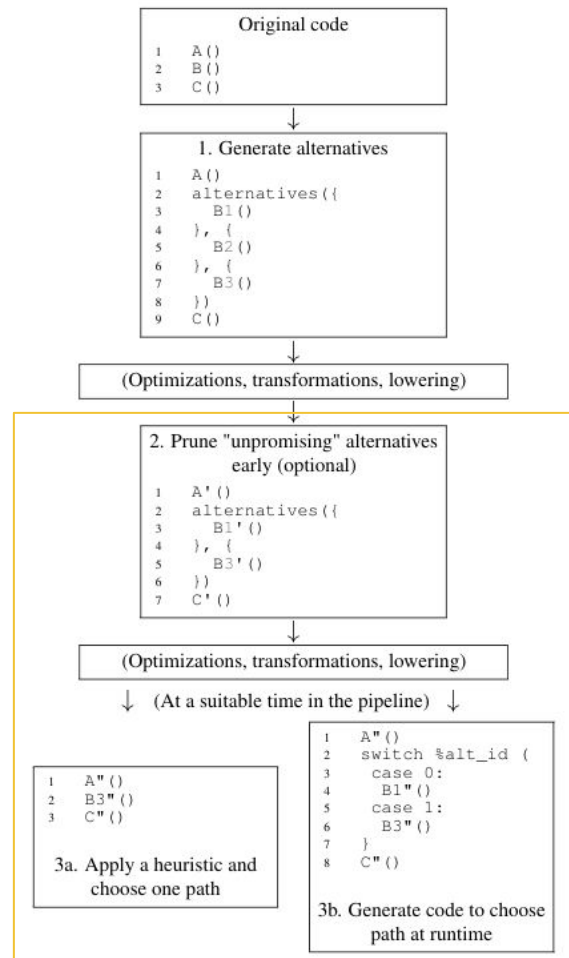  - Discard unpromising alternatives

# Alternative Code Paths

For the left over alternatives:

Easy to use Timing Driven Optimisation workflow

- Compile in *profiling* mode - the compiler generates N versions
- Run program N (or more) times to collect profiling data
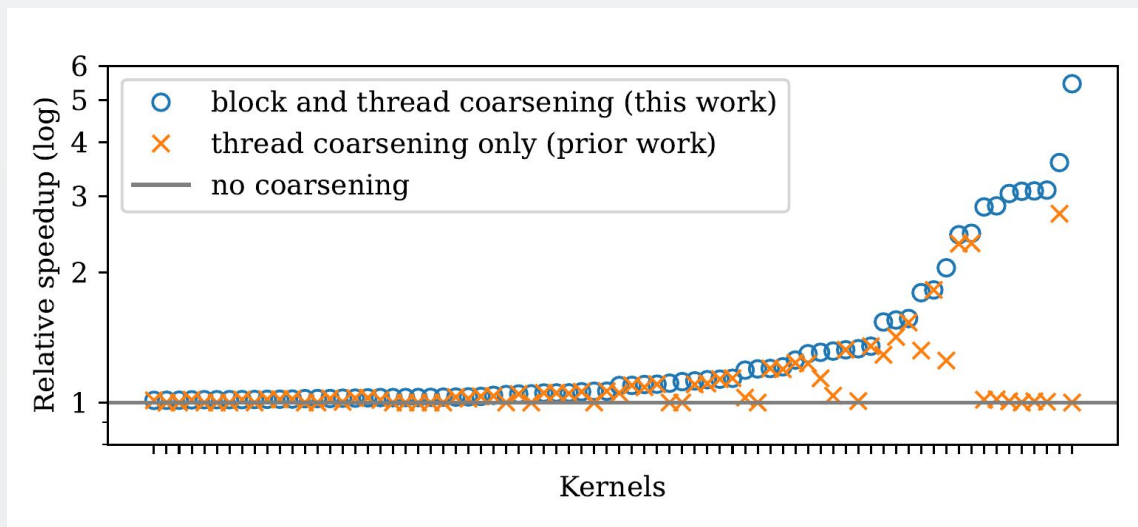- Compile in *optimization* mode - the compiler picks the best configuration for you

# Evaluation

# Evaluation

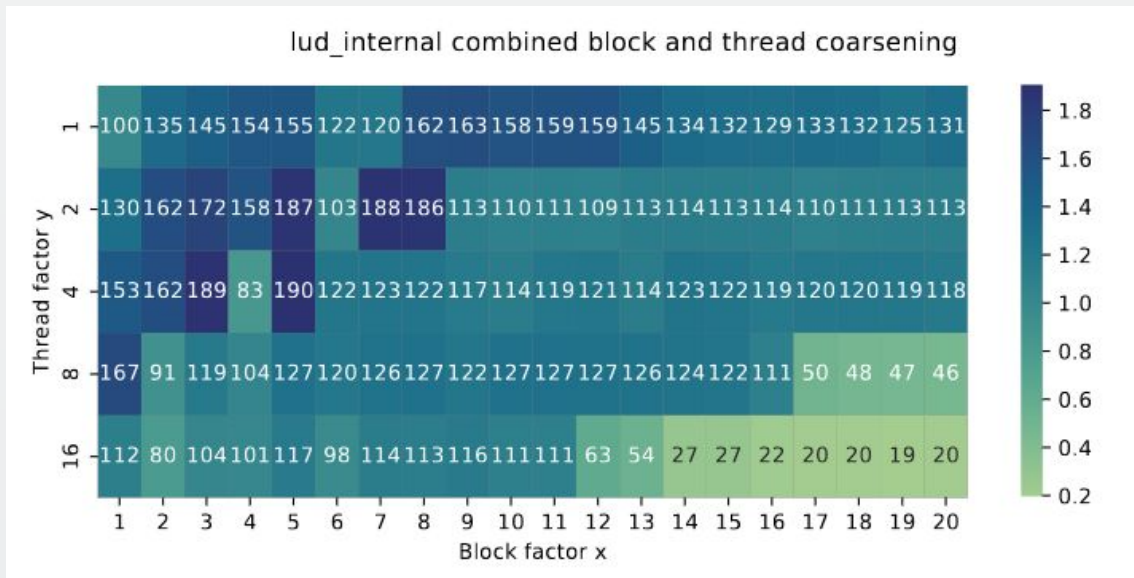On the supported HeCBench and Rodinia CUDA kernels



(best configuration, only the ones with speedup > 1%)

# Some interesting results…

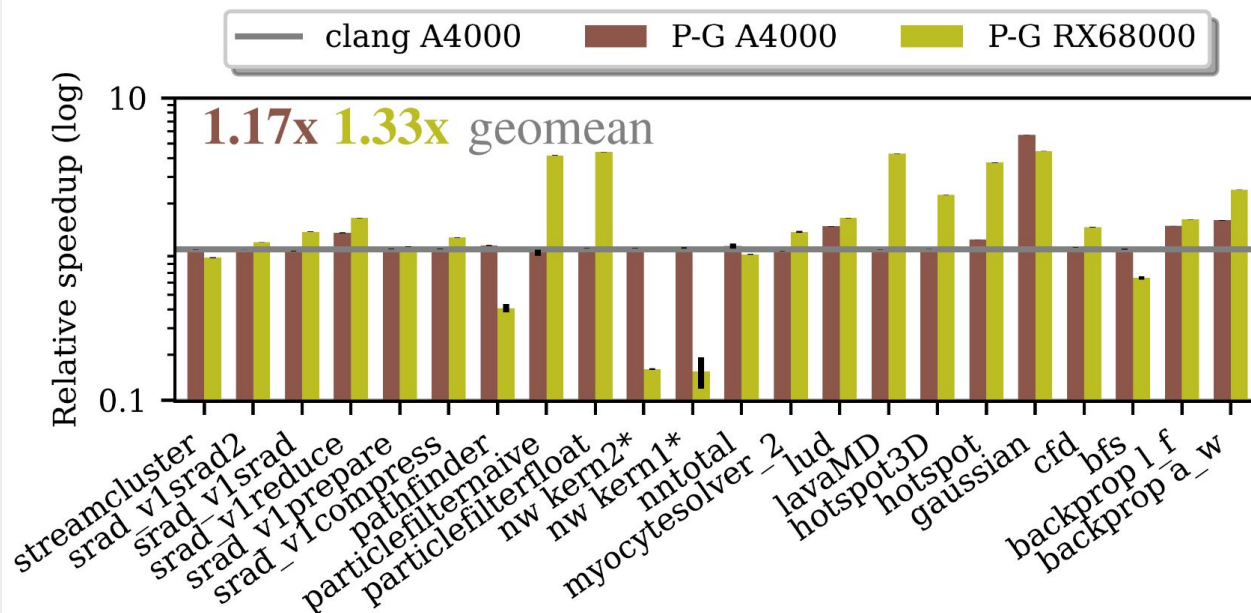Finer granularity of block coarsening needed to maximise performance



Unroll Factor of 5? 3?

# CUDA vs AMD GPU

Rodinia benchmarks

| GPU | Consumer-grade | |
|---|---|---|
| | NVIDIA A4000 | AMD RX6800 |
| Compute Capability | 8.6 | gfx1030 |
| SMs | 48 | 60 |
| FLOPs (f64) | 0.60T | 1.01T |
| FLOPs (f32) | 19.17T | 16.17T |
| Memory Bandwidth | 445 GB/s | 512 GB/s |
| Global Memory | 16 GB | 16 GB |
| L2 Cache | 4 MB | 4MB |
| L1 Cache (Per SM) | 128 KB | 16 KB |

# Conclusion

Performance portability of CUDA code

- Parallel transformations to best utilise available GPU resources
- Timing Driven Optimization framework
- Translation layer to AMD GPU

Nice representation of the computation

⇩

Better (and easier) parallel optimizations